

I. Pen-and-paper

1)

$$\text{basis function } \phi_j(x) = \|w\|_2^j$$

Considerando cada linha como uma observação
 (descolando o output) ficamos com 8 observações

com $D=3$

$$x_1 = [1 \ 10], x_2 = [1 \ 15], x_3 = [0 \ 24], x_4 = [0 \ 23]$$

$$x_5 = [2 \ 0 \ 7], x_6 = [1 \ 1 \ 7], x_7 = [2 \ 0 \ 2], x_8 = [0 \ 2 \ 9]$$

Obtemos uma regressão polinomial com a equação:

$$f(x, w) = \sum_{j=0}^3 w_j \cdot \phi_j(x) = w_0 + w_1 \|x\|_2 + w_2 \|x\|_2^2 + w_3 \|x\|_2^3$$

Para achar os pesos representamos a amostra numa design matrix

$$\Phi = \begin{bmatrix} 1 & \sqrt{2} & \sqrt{2}^2 & \sqrt{2}^3 \\ 1 & \sqrt{5} & (\sqrt{5})^2 & (\sqrt{5})^3 \\ 1 & \sqrt{5} & (\sqrt{5})^2 & (\sqrt{5})^3 \\ 1 & \sqrt{24} & (\sqrt{24})^2 & (\sqrt{24})^3 \\ 1 & \sqrt{53} & (\sqrt{53})^2 & (\sqrt{53})^3 \\ 1 & \sqrt{3} & (\sqrt{3})^2 & (\sqrt{3})^3 \\ 1 & \sqrt{5} & (\sqrt{5})^2 & (\sqrt{5})^3 \\ 1 & \sqrt{85} & (\sqrt{85})^2 & (\sqrt{85})^3 \end{bmatrix}$$

ou que

$$\| [1 \ 10] \| = \sqrt{2}$$

$$\| [1 \ 15] \| = \sqrt{5}$$

$$\| [0 \ 24] \| = \sqrt{24}$$

$$\| [0 \ 23] \| = \sqrt{53}$$

$$\| [2 \ 0 \ 7] \| = \sqrt{5}$$

$$\| [1 \ 1 \ 7] \| = \sqrt{3}$$

$$\| [2 \ 0 \ 2] \| = \sqrt{5}$$

$$\| [0 \ 2 \ 9] \| = \sqrt{85}$$

$$t = \begin{bmatrix} 3 \\ 2 \\ 0 \\ 4 \\ 5 \\ 1 \\ 2 \\ 5 \end{bmatrix}$$

Onde a solução fechada é:

$$w = (\Phi^T \cdot \Phi)^{-1} \cdot \Phi^T \cdot t =$$

$$w = \begin{bmatrix} 4,583521221 \\ -1,687204806 \\ 0,3377373254 \\ -0,07330674257 \end{bmatrix}$$

ficando assim com a equação que representa a regressão polinomial:

$$f(x, w) = 4,583521 - 1,6872048 \|x\|_2 + 0,337737 \|x\|_2^2 - 0,0733067 \|x\|_2^3$$

2)

Consideremos 2 observações de teste
 $x_9 = [2 \ 00]$ e $x_{10} = [1 \ 27]$

Para acharmos a RMSE temos que prever os valores de cada observação no modelo de regressão polinomial para isso achamos o valor de x_9 e x_{10} através da função descreve a regressão polinomial

$$\|x_9\|_2 = \|[2 \ 00]\|_2 = 2$$

$$\|x_{10}\|_2 = \|[1 \ 27]\|_2 = \sqrt{6}$$

$$F(x_9, w) = 4,583521 - 7,6872048 \times 2 + 0,33773773 \times 2^2 - 0,07330674 \times 2^3 = 2,45367$$

$$F(x_{10}, w) = 4,583521 - 7,6872048 \times \sqrt{6} + 0,33773773 \times \sqrt{6}^2 - 0,07330674 \times \sqrt{6}^3 = 2,28759$$

Dado os valores anteriores, e que temos 2 observações sabemos que:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} = \sqrt{\frac{(2,45367 - 2)^2 + (2,28759 - 4)^2}{2}} \approx 1,25672$$

3)

Fazendo uma equal depth-binarization
em y_3 e output, ficamos com:

	y_1	y_2	y_3	output
x_1	1	1	N	N
x_2	1	1	P	N
x_3	0	2	P	N
x_4	1	2	N	N
x_5	2	0	P	P
x_6	1	1	N	P
x_7	2	0	N	P
x_8	0	2	P	P

output $\Rightarrow 0$

$$E_{\text{start}} = E\left(\frac{y_1}{2}, \frac{y_2}{2}\right) = -\left(\frac{4}{8} \log_2\left(\frac{4}{8}\right) + \frac{4}{8} \log_2\left(\frac{4}{8}\right)\right) = 1$$

$$\begin{array}{c|c|c} y_1=0 & y_2=1 & y_2=2 \\ \downarrow & \downarrow & \downarrow \\ \# \{0=N\}=1 & \# \{0=N\}=3 & \# \{0=N\}=0 \\ \# \{0=P\}=1 & \# \{0=P\}=1 & \# \{0=P\}=2 \end{array}$$

$$E(y_1=0) = -\frac{1}{2} \log_2\left(\frac{2}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) = 1$$

$$E(y_1=1) = -\frac{3}{4} \log_2\left(\frac{1}{4}\right) - \frac{1}{4} \log_2\left(\frac{3}{4}\right) \approx 0,8173$$

$$E(y_1=2) = -\frac{2}{2} \log_2\left(\frac{2}{2}\right) = 0$$

$$\begin{aligned} E_{y_1} &= \frac{2}{8} \times E(y_1=0) + \frac{4}{8} \times E(y_1=1) + \frac{2}{8} \times E(y_1=2) = \\ &= \frac{2}{8} \times 1 + \frac{4}{8} \times 0,8173 + \frac{2}{8} \times 0 = 0,65565 \end{aligned}$$

$$\begin{array}{l|l|l}
 y_2=0 & y_2=1 & y_2=2 \\
 \# \{0=N\}=0 & \# \{0=N\}=2 & \# \{0=N\}=2 \\
 \# \{0=P\}=2 & \# \{0=P\}=1 & \# \{0=P\}=1
 \end{array}$$

$$E(y_2=0) = -\frac{2}{2} \log_2\left(\frac{2}{2}\right) = 0$$

$$E(y_2=1) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) \approx 0,9183$$

$$E(y_2=2) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) \approx 0,9183$$

$$E y_2 = \frac{2}{8} \times E(y_2=0) + \frac{3}{8} \times E(y_2=1) + \frac{3}{8} \times E(y_2=2) =$$

$$= \frac{2}{8} \times 0 + \frac{3}{8} \times 0,9183 + \frac{3}{8} \times 0,9183 =$$

$$\approx 0,688725$$

$$I_G(y_2) = E_{\text{start}} - E(y_2) = 1 - 0,688725 = 0,311275$$

$$\begin{array}{l|l}
 y_3=N & y_3=P \\
 \# \{0=N\}=2 & \# \{0=N\}=2 \\
 \# \{0=P\}=2 & \# \{0=P\}=2
 \end{array}$$

$$E(y_3=N) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

$$E(y_3=P) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

$$E(y_3) = \frac{4}{8} E(y_3=N) + \frac{4}{8} E(y_3=P) = \frac{4}{8} \times 2 = 1$$

$$I_G(y_3) = E_{\text{start}} - E(y_3) = 1 - 1 = 0$$

Com isto concluímos que o maior I_G é 0 de y_1 entre y_1 e o nó da árvore e ficamos com:

$$\begin{array}{l|l|l}
 y_1=0 & y_1=1 & y_1=2 \\
 \begin{array}{l} y_2 \ y_3 \ \text{output} \\ 2 \ P \ N \\ 2 \ P \ P \end{array} & \begin{array}{l} y_2 \ y_3 \ \text{output} \\ 1 \ N \ N \\ 1 \ P \ N \\ 2 \ N \ N \\ 1 \ N \ P \end{array} & \begin{array}{l} y_2 \ y_3 \ \text{output} \\ 0 \ P \ P \\ 0 \ N \ P \end{array}
 \end{array}$$

No entanto, as partições ($y_1=0$ e $y_1=1$) ainda têm incerteza, então temos que repetir o processo para ambas

$$E(\text{start } y_1=0) = -\frac{1}{2} \times \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\begin{array}{l|l} y_2=2 & y_3=p \\ \# \{0=N\}=1 & \# \{0=N\}=1 \\ \# \{0=p\}=1 & \# \{0=p\}=1 \end{array}$$

$$E(y_1=0 \wedge y_2=2) = -\frac{1}{2} \log_2 \left(\frac{1}{2} \right)$$

$$I_G(y_1=0 \wedge y_2=2) = 1 - 0,5 = 0,5$$

$$E(y_1=0 \wedge y_2=2 \wedge y_3=p) = -\frac{1}{2} \log_2 \left(\frac{1}{2} \right) = 0,5$$

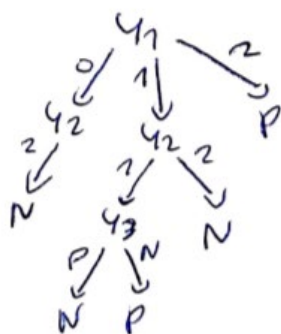
$$I_G(y_1=0 \wedge y_2=2 \wedge y_3=p) = 1 - 0,5 = 0,5$$

Dado que o I_G de ambos é igual escolhemos um qualquer neste caso que se $y_2=2 \wedge y_3=p$ então o output é N .

Também pela tabela acima podemos concluir que se $y_1=1 \wedge y_2=2$ então $O=N$ e que se $y_1=1 \wedge y_2=1 \wedge y_3=p$ então $O=N$.

Pelo caso referido em cima como o I_G se $y_1=1 \wedge y_2=1 \wedge y_3=N$ ser o mesmo escolhe-se que nesse caso o output é P .

Assim sendo, ficamos com a árvore



Aprendizagem 2021/22
Homework II – Group 072

4)

Como podemos ver pela RMSE e a accuracy, o modelo não preve bem os dados de teste devido á RMSE ser alta e a accuracy baixa. Porém, podemos ver que para X9 o valor previsto é muito perto do atual (output), contudo, também podemos ver que o X10 é um noise point e que afeta drásticamente a previsão do modelo, dado que o valor previsto para X10 (2,28159) é muito menor que o valor do output (4).

$$\text{accuracy} = \frac{\text{n}^{\circ} \text{previsões corretas}}{\text{total número de previsões}} = \frac{1}{2} = 0,5$$

II. Programming and critical analysis

5) Answer 5

6) Answer 6

7) Answer 7

III. APPENDIX

END